

Regression: want $MSE[f] = E[(Y - f(x))^2]$
to be small

• If optimal regression prediction is
 $\mu(x) = E(Y|X=x)$, but

$f(y|x)$ is unknown ($\mu(x) = \int_{-\infty}^{\infty} y f(y|x) dy$)

• Can't just do $E(Y|X=x)$ is $\frac{1}{\#\{i: X_i = x\}} \sum_{i: X_i = x} Y_i$
since ~~over~~ that value of x might not appear
in the data

• Will use linear smoother so

$$\hat{\mu}(x) = \sum_{i=1}^n w(x_i, x) Y_i$$

(x with no subscript is the point we are trying to predict)
weight of Y_i depends on what x_i is and on where we're predicting

• common weight functions:

□ global mean $w(x_i, x) = \frac{1}{n}$ if x_i is closer to x than any other x_i

□ nearest neighbors $w(x_i, x) = \begin{cases} 1 & \text{if } x_i \text{ is closest} \\ 0 & \text{otherwise} \end{cases}$

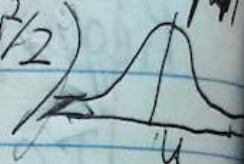
□ k nearest neighbors $w(x_i, x) = \begin{cases} \frac{1}{k} & \text{if } x_i \text{ is one of the } k \text{ closest} \\ 0 & \text{otherwise} \end{cases}$

• increasing k reduces variance but increases bias

• for knn : $\hat{\mu}(x) = \frac{1}{k} \sum_{i: i \in KNN(x)} \mu(x_i) + \frac{1}{k} \sum_{i: i \in KNN(x)} Y_i$

$$\square w(x_i, x) = \frac{k\left(\frac{x_i - x}{h}\right)}{\sum_{j=1}^n k\left(\frac{x_j - x}{h}\right)}$$

where $k(u)$ is a smooth pdf
(say $k(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}$)



h is scale/bandwidth

- data points with $|x_i - x| < h$ matter a lot
- data points with $|x_i - x| > h$ hardly matter

The function k is called kernel and this is kernel regression

- $h \rightarrow 0$ looks like nearest neighbors (less bias, more var)
- $h \rightarrow \infty$ estimate flattens out to global mean
- $\sum_{i=1}^n w(x_i, x, h) = 1$ for kernel regression
- These are nonparametric (global mean, k-NN, kernel)

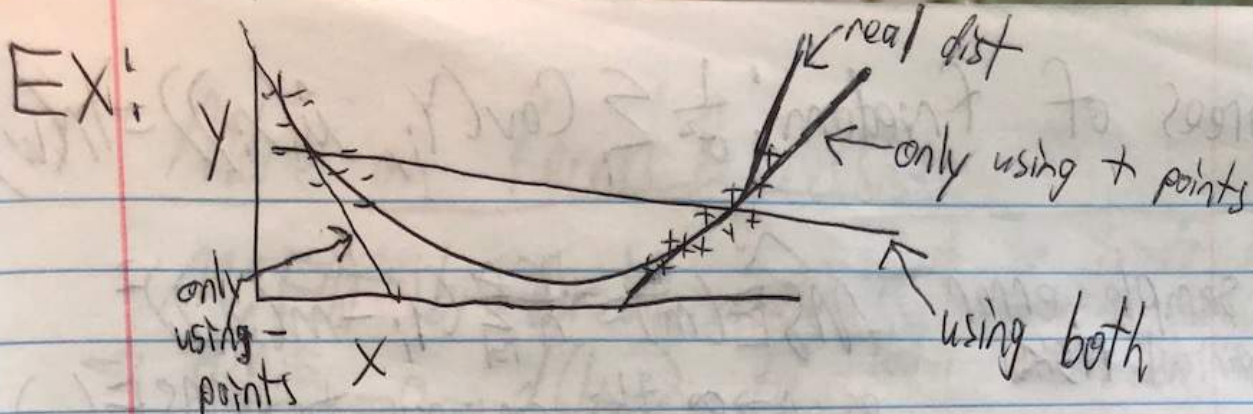
Linear regression

A new predicted y is $x_{\text{new}} \hat{\beta} = x_{\text{new}} (x^T x)^{-1} x^T y$

$$\Rightarrow x_{\text{new}} \hat{\beta} = \sum_{i=1}^n x_{\text{new}} \frac{x_i y_i}{n \sigma_x^2}$$

- (i) This doesn't depend on x we give more weight to ~~small~~ points closer to x_{new} and
- (ii) $w(x_i, x)$ goes up as x_i moves away from \bar{x}
- (iii) this only makes sense if μ is really linear

• If the true $\mu(x)$ is not linear then there is always an optimal linear β , but it changes with the distribution of x . Even though best linear model it could still be bad.



Optimal linear predictor of Y from X : $(\alpha, \beta) = \underset{a, b}{\operatorname{argmin}} E[(Y - (a + bX))^2]$
 $\Rightarrow \alpha = E[Y] - \beta E[X]$ and $\beta = \frac{\operatorname{Cov}(X, Y)}{\operatorname{Var}(X)}$
 which says $\operatorname{Cov}(X, Y) = \beta \operatorname{Var}(X)$

$$\hat{\mu} = WY = \begin{bmatrix} w(x_1, x_1) & \dots & w(x_1, x_n) \\ \vdots & \ddots & \vdots \\ w(x_n, x_1) & \dots & w(x_n, x_n) \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\hat{\mu} = W\mu + W\varepsilon$$

~~$\hat{\mu} = W\mu + W\varepsilon$~~

$$\Rightarrow \operatorname{bias}(\hat{\mu}) = \mu - W\mu = (I - W)\mu$$

$$\Rightarrow \operatorname{Var}(\hat{\mu}) = \operatorname{Var}(W\varepsilon) = \sigma^2 W W^T$$

assuming that $\operatorname{Var}(\varepsilon_i) = \operatorname{Var}(\varepsilon_j) = 0 \quad \forall i, j$ and
 $\operatorname{Cov}(\varepsilon_i, \varepsilon_j) = 0$ (unless $i = j$)

• for linear models $W = X(X^T X)^{-1} X^T = H$ (hat matrix)

degrees of freedom: $\frac{1}{\sigma^2} \sum_i \text{Cov}(Y_i, \hat{u}(x_i)) = \text{tr}(K)$

In sample error (empirical ~~MSE~~) $\hat{MSE}(m) = \frac{1}{n} \sum_{i=1}^n (Y_i - m(x_i))^2$
as $n \rightarrow \infty$ this converges to $MSE(m)$

- Usually we have \hat{m} estimated from the data (so it's not fixed in advance)
- Minimizing \hat{MSE} partially reduces MSE but also partially memorizes noise

How to avoid memorizing noise?

- a) use theory to estimate/bound $MSE(\hat{m}) - \hat{MSE}(\hat{m})$ (degrees of freedom adjustments)
- b) Use new data to evaluate different models (calculate \hat{MSE} for each candidate model with hold-out/validation data)
 - Can't use the new data to fit model
- c) Fake having new data with cross-validation

K-folds cross validation:

- Randomly divide data into k equal parts
- each part gets used once as the testing set and use the other $k-1$ folds as the training set
- Leave one out cross validation (LOOCV) is using n -folds (only 1 point in testing set)
- After using ~~fitting~~ k -folds to pick the model ~~and~~ fit the chosen model to all of the data